

RATES OF CONVERGENCE IN ACTIVE LEARNING

BY STEVE HANNEKE¹

Carnegie Mellon University

We study the rates of convergence in generalization error achievable by active learning under various types of label noise. Additionally, we study the general problem of model selection for active learning with a nested hierarchy of hypothesis classes and propose an algorithm whose error rate provably converges to the best achievable error among classifiers in the hierarchy at a rate adaptive to both the complexity of the optimal classifier and the noise conditions. In particular, we state sufficient conditions for these rates to be dramatically faster than those achievable by passive learning.

1. Introduction. *Active learning* refers to a family of powerful supervised learning protocols capable of producing more accurate classifiers while using a smaller number of labeled data points than traditional (passive) learning methods. Here we study a variant known as *pool-based* active learning, in which a learning algorithm is given access to a large pool of unlabeled data (i.e., only the covariates are visible), and is allowed to sequentially request the label (response variable) of any particular data points from that pool. The objective is to learn a function that accurately predicts the labels of new points, while minimizing the number of label requests. Thus, this is a type of sequential design scenario for a function estimation problem. This contrasts with passive learning, where the labeled data are sampled at random. In comparison, by more carefully selecting which points should be labeled, active learning can often significantly decrease the total amount of effort required for data annotation. This can be particularly interesting for tasks where unlabeled data are available in abundance, but label information comes only through significant effort or cost.

Received August 2009; revised June 2010.

¹Supported by the NSF Grant IIS-0713379 and an IBM Ph.D. Fellowship.

AMS 2000 subject classifications. Primary 62L05, 68Q32, 62H30, 68T05; secondary 68T10, 68Q10, 68Q25, 68W40, 62G99.

Key words and phrases. Active learning, sequential design, selective sampling, statistical learning theory, oracle inequalities, model selection, classification.

This is an electronic reprint of the original article published by the Institute of Mathematical Statistics in *The Annals of Statistics*, 2011, Vol. 39, No. 1, 333–361. This reprint differs from the original in pagination and typographic detail.

Recently, there have been a series of exciting advances on the topic of active learning with arbitrary classification noise (the so-called *agnostic* PAC model [22]), resulting in several new algorithms capable of achieving improved convergence rates compared to passive learning under certain conditions. The first, proposed by Balcan, Beygelzimer and Langford [6] was the A^2 (agnostic active) algorithm, which provably never has significantly worse rates of convergence than passive learning by empirical risk minimization. This algorithm was later analyzed in detail in [19], where it was found that a complexity measure called the *disagreement coefficient* characterizes the worst-case convergence rates achieved by A^2 for any given hypothesis class, data distribution and best achievable error rate in the class. The next major advance was by Dasgupta, Hsu and Monteleoni [14], who proposed a new algorithm, and proved that it improves the dependence of the convergence rates on the disagreement coefficient compared to A^2 . Both algorithms are defined below in Section 3. While all of these advances are encouraging, they are limited in two ways. First, the convergence rates that have been proven for these algorithms typically only improve the dependence on the magnitude of the noise (more precisely, the noise rate of the hypothesis class), compared to passive learning. Thus, in an asymptotic sense, for nonzero noise rates these results represent at best a constant factor improvement over passive learning. Second, these results are limited to learning with a fixed hypothesis class of limited expressiveness, so that convergence to the Bayes error rate is not always a possibility.

On the first of these limitations, recent work by Castro and Nowak [12] on learning threshold classifiers discovered that if certain parameters of the noise distribution are *known* (namely, parameters related to Tsybakov’s margin conditions), then we can achieve strict improvements in the asymptotic convergence rate via a specific active learning algorithm designed to take advantage of that knowledge for thresholds. Subsequently, Balcan, Broder and Zhang [7] proved a similar result for linear separators in higher dimensions, and Castro and Nowak [12] showed related improvements for the space of boundary fragment classes (under a somewhat stronger assumption than Tsybakov’s). However, these works left open the question of whether such improvements could be achieved by an algorithm that does not explicitly depend on the noise conditions (i.e., in the *agnostic* setting), and whether this type of improvement is achievable for more general families of hypothesis classes, under the usual complexity restrictions (e.g., VC class, entropy conditions, etc.). In a personal communication, John Langford and Rui Castro claimed A^2 achieves these improvements for the special case of threshold classifiers (a special case of this also appeared in [9]). However, there remained an open question of whether such rate improvements could be generalized to hold for arbitrary hypothesis classes. In Section 4, we provide this generalization. We analyze the rates achieved by A^2 under Tsybakov’s

noise conditions [26, 28]; in particular, we find that these rates are strictly superior to the known rates for passive learning, when the disagreement coefficient is finite. We also study a novel modification of the algorithm of Dasgupta, Hsu and Monteleoni [14], proving that it improves upon the rates of A^2 in its dependence on the disagreement coefficient.

Additionally, in Section 5, we address the second limitation by proposing a general model selection procedure for active learning with an arbitrary structure of nested hypothesis classes. If the classes have restricted expressiveness (e.g., VC classes), the error rate for this algorithm converges to the best achievable error by any classifier in the structure, at a rate that adapts to the noise conditions and complexity of the optimal classifier. In general, if the structure is constructed to include arbitrarily good approximations to any classifier, the error converges to the Bayes error rate in the limit. In particular, if the Bayes optimal classifier is in some class within the structure, the algorithm performs nearly as well as running an agnostic active learning algorithm on that single hypothesis class, thus preserving the convergence rate improvements achievable for that class.

2. Definitions and notation. In the active learning setting, there is an *instance space* \mathcal{X} , a *label space* $\mathcal{Y} = \{-1, +1\}$ and some fixed distribution \mathcal{D}_{XY} over $\mathcal{X} \times \mathcal{Y}$, with marginal \mathcal{D}_X over \mathcal{X} . The restriction to binary classification ($\mathcal{Y} = \{-1, +1\}$) is intended to simplify the discussion; however, everything below generalizes quite naturally to multiclass classification (where $\mathcal{Y} = \{1, 2, \dots, k\}$).

There are two sequences of random variables: X_1, X_2, \dots and Y_1, Y_2, \dots , where each (X_i, Y_i) pair is independent of the others, and has joint distribution \mathcal{D}_{XY} . However, the learning algorithm is only permitted direct access to the X_i values (unlabeled data points), and must request the Y_i values one at a time, sequentially. That is, the algorithm picks some index i to observe the Y_i value, then after observing it, picks another index i' to observe the $Y_{i'}$ label value, etc. We are interested in studying the rate of convergence of the error rate of the classifier output by the learning algorithm, in terms of the number of label requests it has made. To simplify the discussion, we will think of the data sequence as being essentially inexhaustible, and will study $(1 - \delta)$ -confidence bounds on the error rate of the classifier produced by an algorithm permitted to make at most n label requests, for a fixed value $\delta \in (0, 1/2)$. The actual number of (unlabeled) data points the algorithm uses will be made clear in the proofs (typically close to the number of points needed by passive learning to achieve the stated error guarantee).

A *hypothesis class* \mathbb{C} is any set of measurable classifiers $h: \mathcal{X} \rightarrow \mathcal{Y}$. We will denote by d the VC dimension of \mathbb{C} (see, e.g., [11, 15, 30–32]). For any measurable $h: \mathcal{X} \rightarrow \mathcal{Y}$ and distribution \mathcal{D} over $\mathcal{X} \times \mathcal{Y}$, define the *error rate* of h as $er_{\mathcal{D}}(h) = \mathbb{P}_{(X,Y) \sim \mathcal{D}}\{h(X) \neq Y\}$; when $\mathcal{D} = \mathcal{D}_{XY}$, we abbreviate

this as $er(h)$. This simply represents the risk under the 0–1 loss. We also define the *conditional error rate*, given a set $R \subseteq \mathcal{X}$, as $er(h|R) = \mathbb{P}\{h(X) \neq Y|X \in R\}$. Let $\nu = \inf_{h \in \mathbb{C}} er(h)$, called the *noise rate* of \mathbb{C} . For any $x \in \mathcal{X}$, let $\eta(x) = \mathbb{P}\{Y = 1|X = x\}$, let $h^*(x) = 2\mathbb{1}[\eta(x) \geq 1/2] - 1$ and let $\nu^* = er(h^*)$. We call h^* the *Bayes optimal classifier* and ν^* the *Bayes error rate*. Additionally, define the *diameter* of any set of classifiers V as $\text{diam}(V) = \sup_{h_1, h_2 \in V} \mathbb{P}\{h_1(X) \neq h_2(X)\}$, and for any $\varepsilon > 0$, define the diameter of the ε -*minimal set* of V as $\text{diam}(\varepsilon; V) = \text{diam}(\{h \in V : er(h) - \inf_{h' \in V} er(h') \leq \varepsilon\})$.

For a classifier h , and a sequence $S = \{(x_1, y_1), (x_2, y_2), \dots, (x_m, y_m)\} \in (\mathcal{X} \times \mathcal{Y})^m$, let $er_S(h) = \frac{1}{|S|} \sum_{(x,y) \in S} \mathbb{1}[h(x) \neq y]$ denote the *empirical error rate* on S , [and define $er_{\emptyset}(h) = 0$ by convention]. It will often be convenient to make use of sets of (index, label) pairs, where the index is used to uniquely refer to an element of the $\{X_i\}$ sequence (while conveniently also keeping track of relative ordering information); in such contexts, we will overload notation as follows. For a classifier h , and a finite set of (index, label) pairs $S = \{(i_1, y_1), (i_2, y_2), \dots, (i_m, y_m)\} \subset \mathbb{N} \times \mathcal{Y}$, let $er_S(h) = \frac{1}{|S|} \sum_{(i,y) \in S} \mathbb{1}[h(X_i) \neq y]$, (and $er_{\emptyset}(h) = 0$, as before). Thus, $er_S(h) = er_{S'}(h)$, where $S' = \{(X_i, y)\}_{(i,y) \in S}$. For the indexed *true* label sequence, $\mathcal{Z}^{(m)} = \{(1, Y_1), (2, Y_2), \dots, (m, Y_m)\}$, we abbreviate this $er_m(h) = er_{\mathcal{Z}^{(m)}}(h)$, the empirical error on the first m data points.

In addition to the independent interest of understanding the rates achievable here, another primary interest in this setting is to quantify the achievable *improvements*, compared to *passive learning*. In this context, a passive learning algorithm can be formally defined as a function mapping the sequence $\{(X_1, Y_1), (X_2, Y_2), \dots, (X_n, Y_n)\}$ to a classifier \hat{h}_n ; for instance, perhaps the most widely studied family of passive learning methods is that of *empirical risk minimization* (e.g., [23, 27, 30, 31]), which return a classifier $\hat{h}_n \in \arg \min_{h \in \mathbb{C}} er_n(h)$. For the purpose of this comparison, we review known results on passive learning in several contexts below.

2.1. Tsybakov’s noise conditions. Here we describe a particular parametrization of noise distributions, relative to a hypothesis class, often referred to as Tsybakov’s noise conditions [26, 28], or margin conditions. These noise conditions have recently received substantial attention in the passive learning literature, as they describe situations in which the asymptotic minimax convergence rate of passive learning is faster than the worst case $n^{-1/2}$ rate (e.g., [23, 26–28]).

CONDITION 1. There exist finite constants $\mu > 0$ and $\kappa \geq 1$, s.t. $\forall \varepsilon > 0$, $\text{diam}(\varepsilon; \mathbb{C}) \leq \mu \varepsilon^{1/\kappa}$.

This condition is satisfied when, for example,

$$\exists \mu' > 0, \kappa \geq 1 \text{ s.t. } \exists h \in \mathbb{C} : \forall h' \in \mathbb{C} \quad er(h') - \nu \geq \mu' \mathbb{P}\{h(X) \neq h'(X)\}^\kappa,$$

[23]. It is also satisfied when the Bayes optimal classifier is in \mathbb{C} and

$$\exists \mu'' > 0, \alpha \in (0, \infty) \text{ s.t. } \forall \varepsilon > 0 \quad \mathbb{P}\{|\eta(X) - 1/2| \leq \varepsilon\} \leq \mu'' \varepsilon^\alpha,$$

where κ and μ are functions of α and μ'' [26, 28]; in particular, $\kappa = (1 + \alpha)/\alpha$. As we will see, the case where $\kappa = 1$ is particularly interesting; for instance, this is the case when $h^* \in \mathbb{C}$ and $\mathbb{P}\{|\eta(X) - 1/2| > c\} = 1$ for some constant $c \in (0, 1/2)$. Informally, in many cases Condition 1 can be realized in terms of the relation between magnitude of noise and distance to the optimal decision boundary; that is, since in practice the amount of noise in a data point's label is often inversely related to the distance from the decision boundary, a small κ value may often result from having low density near the decision boundary (i.e., large margin); when this is not the case, the value of κ is often determined by how quickly $\eta(x)$ changes as x approaches the decision boundary. See [7, 12, 23, 26–28] for further interpretations of this condition.

It is known that when this condition is satisfied for some $\kappa \geq 1$ and $\mu > 0$, the passive learning method of empirical risk minimization achieves a convergence rate guarantee, holding with probability $\geq 1 - \delta$, of

$$er\left(\arg \min_{h \in \mathbb{C}} er_n(h)\right) - \nu \leq c \left(\frac{d \log n + \log(1/\delta)}{n} \right)^{\kappa/(2\kappa-1)},$$

where c is a (κ and μ -dependent) constant (this follows from [23, 27]; see Appendix B of the supplementary material [20], especially (17) and Lemma 5, for the details). Furthermore, for some hypothesis classes, this is known to be a tight bound (up to the log factor) on the minimax convergence rate, so that there is *no* passive learning algorithm for these classes for which we can guarantee a faster convergence rate, given that the guarantee depends on \mathcal{D}_{XY} only through μ and κ [12, 28] (see also Appendix D of the supplementary material [20]).

2.2. Disagreement coefficient. The disagreement coefficient, introduced in [19], is a measure of the complexity of an active learning problem, which has proven quite useful for analyzing the convergence rates of certain types of active learning algorithms: for example, the algorithms of [6, 10, 13, 14]. Informally, it quantifies how much disagreement there is among a set of classifiers relative to how close to some h they are. The following is a version of its definition, which we will use extensively below. For any hypothesis class \mathbb{C} and $V \subseteq \mathbb{C}$, let

$$\text{DIS}(V) = \{x \in \mathcal{X} : \exists h_1, h_2 \in V \text{ s.t. } h_1(x) \neq h_2(x)\}.$$

For $r \in [0, 1]$ and measurable $h: \mathcal{X} \rightarrow \mathcal{Y}$, let

$$B(h, r) = \{h' \in \mathbb{C} : \mathbb{P}\{h(X) \neq h'(X)\} \leq r\}.$$

DEFINITION 1. The *disagreement coefficient* of h with respect to \mathbb{C} under \mathcal{D}_X is defined as

$$\theta_h = \sup_{r > r_0} \frac{\mathbb{P}(\text{DIS}(B(h, r)))}{r},$$

where $r_0 = 0$ (though see Appendix A.1 for alternative possibilities for r_0).

DEFINITION 2. We further define the disagreement coefficient for the hypothesis class \mathbb{C} with respect to the target distribution \mathcal{D}_{XY} as $\theta = \liminf_{k \rightarrow \infty} \theta_{h^{[k]}}$, where $\{h^{[k]}\}$ is any sequence in \mathbb{C} with $er(h^{[k]})$ monotonically decreasing to ν ; [by convention, take every $h^{[k]} \in \arg \min_{h \in \mathbb{C}} er(h)$ if the minimum is achieved].

In Definition 1, it is conceivable that $\text{DIS}(B(h, r))$ may sometimes not be measurable. In such cases, we can define $\mathbb{P}(\text{DIS}(B(h, r)))$ as the *outer* measure [29], so that it remains well defined. We continue this practice below, letting \mathbb{P} and \mathbb{E} (and indeed any reference to “probability”) refer to the outer expectation and measure in any context for which this is necessary.

Because of its simple intuitive interpretation, measuring the amount of disagreement in a local neighborhood of some classifier h , the disagreement coefficient has the wonderful property of being relatively simple to calculate for a wide range of learning problems, especially when those problems have a natural geometric representation. To illustrate this, we will go through a few simple examples from [19].

Consider the hypothesis class of thresholds h_z on the interval $[0, 1]$ [for $z \in (0, 1)$], where $h_z(x) = +1$ iff $x \geq z$. Furthermore, suppose \mathcal{D}_X is uniform on $[0, 1]$. In this case, it is clear that the disagreement coefficient is 2, since for sufficiently small r , the region of disagreement of $B(h_z, r)$ is $[z - r, z + r]$, which has probability mass $2r$. In other words, since the disagreement region grows with r in two disjoint directions, each at rate 1, we have $\theta_{h_z} = 2$.

As a second example, consider the disagreement coefficient for *intervals* on $[0, 1]$. As before, let $\mathcal{X} = [0, 1]$ and \mathcal{D}_X be uniform, but this time \mathbb{C} is the set of intervals $h_{[a, b]}$ such that for $x \in [0, 1]$, $h_{[a, b]}(x) = +1$ iff $x \in [a, b]$ (for $0 < a < b < 1$). In contrast to thresholds, the disagreement coefficients $\theta_{h_{[a, b]}}$ for the space of intervals vary widely depending on the particular $h_{[a, b]}$. Specifically, we have $\theta_{h_{[a, b]}} = \max\{\frac{1}{b-a}, 4\}$. To see this, note that when $0 < r < b - a$, every

interval in $B(h_{[a,b]}, r)$ has its lower and upper boundaries within r of a and b , respectively; thus, $\mathbb{P}(\text{DIS}(B(h_{[a,b]}, r))) \leq 4r$, with equality for sufficiently small r . However, when $r > b - a$, every interval of width $\leq r - (b - a)$ is in $B(h_{[a,b]}, r)$, so that $\mathbb{P}(\text{DIS}(B(h_{[a,b]}, r))) = 1$.

As a slightly more involved example, [19] studies the scenario where \mathcal{X} is the surface of the origin-centered unit sphere in \mathbb{R}^d for $d > 2$, \mathbb{C} is the space of all linear separators whose decision surface passes through the origin, and \mathcal{D}_X is the uniform distribution on \mathcal{X} ; in this case, it turns out $\forall h \in \mathbb{C}$ the disagreement coefficient θ_h satisfies

$$\frac{\pi}{4}\sqrt{d} \leq \theta_h \leq \pi\sqrt{d}.$$

The disagreement coefficient has many interesting properties that can help to bound its value for a given hypothesis class and distribution. We list a few elementary properties below. Their proofs, which are quite short and follow directly from the definition, are left as easy exercises.

LEMMA 1 (Close marginals [19]). *Suppose $\exists \lambda \in (0, 1]$ s.t. for any measurable set $A \subseteq \mathcal{X}$, $\lambda \mathbb{P}_{\mathcal{D}_X}(A) \leq \mathbb{P}_{\mathcal{D}'_X}(A) \leq \frac{1}{\lambda} \mathbb{P}_{\mathcal{D}_X}(A)$. Let $h: \mathcal{X} \rightarrow \mathcal{Y}$ be a measurable classifier, and suppose θ_h and θ'_h are the disagreement coefficients for h with respect to \mathbb{C} under \mathcal{D}_X and \mathcal{D}'_X , respectively. Then*

$$\lambda^2 \theta_h \leq \theta'_h \leq \frac{1}{\lambda^2} \theta_h.$$

LEMMA 2 (Finite mixtures). *Suppose $\exists \alpha \in [0, 1]$ s.t. for any measurable set $A \subseteq \mathcal{X}$, $\mathbb{P}_{\mathcal{D}_X}(A) = \alpha \mathbb{P}_{\mathcal{D}_1}(A) + (1 - \alpha) \mathbb{P}_{\mathcal{D}_2}(A)$. For a measurable $h: \mathcal{X} \rightarrow \mathcal{Y}$, let $\theta_h^{(1)}$ be the disagreement coefficient with respect to \mathbb{C} under \mathcal{D}_1 , $\theta_h^{(2)}$ be the disagreement coefficient with respect to \mathbb{C} under \mathcal{D}_2 , and θ_h be the disagreement coefficient with respect to \mathbb{C} under \mathcal{D}_X . Then*

$$\theta_h \leq \theta_h^{(1)} + \theta_h^{(2)}.$$

LEMMA 3 (Finite unions). *Suppose $h \in \mathbb{C}_1 \cap \mathbb{C}_2$ is a classifier s.t. the disagreement coefficient with respect to \mathbb{C}_1 under \mathcal{D}_X is $\theta_h^{(1)}$ and with respect to \mathbb{C}_2 under \mathcal{D}_X is $\theta_h^{(2)}$. Then if θ_h is the disagreement coefficient with respect to $\mathbb{C} = \mathbb{C}_1 \cup \mathbb{C}_2$ under \mathcal{D}_X , we have that*

$$\max\{\theta_h^{(1)}, \theta_h^{(2)}\} \leq \theta_h \leq \theta_h^{(1)} + \theta_h^{(2)}.$$

In fact, even if $h \notin \mathbb{C}_1 \cap \mathbb{C}_2$, we still have $\theta_h \leq \theta_h^{(1)} + \theta_h^{(2)} + 2$.

See [8, 10, 14, 16, 19, 33] for further discussions of various uses of the disagreement coefficient and related notions and extensions in active learning. In particular, Friedman [16] proves that any hypothesis class and distribution satisfying certain general regularity conditions will admit finite constant bounds on θ . Also, Wang [33] bounds the disagreement coefficient for certain nonparametric hypothesis classes, characterized by smoothness of their decision surfaces. Additionally, Beygelzimer, Dasgupta and Langford [10] present an interesting analysis using a natural extension of the disagreement coefficient to study active learning with a larger family of loss functions beyond 0–1 loss.

The disagreement coefficient has deep connections to several other quantities, such as doubling dimension [25] and VC dimension [30]. Additionally, a related quantity, referred to as the “capacity function,” was studied in the 1980s by Alexander in the passive learning literature, in the context of ratio-type empirical processes [2–4] and recently was further developed by Giné and Koltchinskii [17]; interestingly, in this latter work, Giné and Koltchinskii study a localized version of the capacity function, which in our present context can essentially be viewed as the function $\tau(r) = \mathbb{P}(\text{DIS}(B(h, r))) / r$, so that $\theta_h = \sup_{r > r_0} \tau(r)$.

3. General algorithms. We begin the discussion of the algorithms we will analyze by noting the underlying inspiration that unifies them. Specifically, at this writing, all of the published general-purpose agnostic active learning algorithms achieving nontrivial improvements are derivatives of a basic technique proposed by Cohn, Atlas and Ladner [13] for the realizable active learning problem. Under the assumption that there exists a perfect classifier in \mathbb{C} , they proposed an algorithm which processes unlabeled data points in sequence, and for each one it determines whether there is a classifier in \mathbb{C} consistent with all previously observed labels that predicts +1 for this new point *and* one that predicts –1 for this new point; if so, the algorithm requests the label, and otherwise it does not request the label; after n label requests, the algorithm returns any classifier consistent with all observed labels. In some sense, this algorithm corresponds to the very least we could expect of an active learning algorithm, as it never requests the label of a point it can derive from known information, but otherwise makes no effort to search for informative data points. The idea is appealing, not only for its simplicity, but also for its extremely efficient use of unlabeled data; in fact, under the stated assumption, the algorithm produces a classifier consistent with the labels of *all* of the unlabeled data it processes, including those it does *not* request the labels of.

We can equivalently think of this algorithm as maintaining two sets: $V \subseteq \mathbb{C}$ is the set of candidate hypotheses still under consideration, and $R = \text{DIS}(V)$ is their region of disagreement. We can then think of the algorithm

as requesting a random labeled point from the conditional distribution of \mathcal{D}_{XY} given that $X \in R$, and subsequently removing from V any classifier inconsistent with the observed label. A formal definition of the algorithm is given as follows.

Algorithm 0

Input: hypothesis class \mathbb{C} , label budget n

Output: classifier $\hat{h}_n \in \mathbb{C}$

-
0. $V_0 \leftarrow \mathbb{C}$, $t \leftarrow 0$
 1. For $m = 1, 2, \dots$
 2. If $X_m \in \text{DIS}(V_t)$,
 3. Request Y_m
 4. $t \leftarrow t + 1$
 5. $V_t \leftarrow \{h \in V_{t-1} : h(X_m) = Y_m\}$
 6. If $t = n$ or $\{m' > m : X_{m'} \in \text{DIS}(V_t)\} = \emptyset$, Return any $\hat{h}_n \in V_t$
-

The algorithms described below for the problem of active learning with label noise each represent noise-robust variants of this basic idea. They work to reduce the set of candidate hypotheses, while only requesting the labels of points in the region of disagreement of these candidates. The trick is to only remove a classifier from the candidate set once we have high statistical confidence that it is worse than some other candidate classifier so that we never remove the best classifier. However, the two algorithms differ somewhat in the details of how that confidence is calculated.

3.1. Algorithm 1. The first noise-robust algorithm we study, originally proposed by Balcan, Beygelzimer and Langford [6], is typically referred to as A^2 for *Agnostic Active*. This was historically the first general-purpose agnostic active learning algorithm shown to achieve improved error guarantees for certain learning problems in certain ranges of n and ν . Below is a variant of this algorithm. It is defined in terms of two functions: UB and LB . These represent upper and lower confidence bounds on the error rate of a classifier from \mathbb{C} with respect to an arbitrary sampling distribution, as a function of a labeled sequence sampled according to that distribution. Some steps in the algorithm require calculating certain probabilities, such as $\mathbb{P}(\text{DIS}(V))$ or $\mathbb{P}(R)$; later, we discuss replacing these with appropriate estimators.

Algorithm 1

Input: hypothesis class \mathbb{C} , label budget n , confidence δ , functions UB and LB

Output: classifier \hat{h}_n

-
0. $V \leftarrow \mathbb{C}$, $R \leftarrow \text{DIS}(\mathbb{C})$, $Q \leftarrow \emptyset$, $m \leftarrow 0$
 1. For $t = 1, 2, \dots, n$
 2. If $\mathbb{P}(\text{DIS}(V)) \leq \frac{1}{2}\mathbb{P}(R)$
 3. $R \leftarrow \text{DIS}(V)$; $Q \leftarrow \emptyset$
 4. If $\mathbb{P}(R) \leq 2^{-n}$, Return any $\hat{h}_n \in V$
 5. $m \leftarrow \min\{m' > m : X_{m'} \in R\}$
 6. Request Y_m and let $Q \leftarrow Q \cup \{(m, Y_m)\}$
 7. $V \leftarrow \{h \in V : LB(h, Q, \delta/n) \leq \min_{h' \in V} UB(h', Q, \delta/n)\}$
 8. $h_t \leftarrow \arg \min_{h \in V} UB(h, Q, \delta/n)$
 9. $\beta_t \leftarrow (UB(h_t, Q, \delta/n) - \min_{h \in V} LB(h, Q, \delta/n))\mathbb{P}(R)$
 10. Return $\hat{h}_n = h_{\hat{t}}$, where $\hat{t} = \arg \min_{t \in \{1, 2, \dots, n\}} \beta_t$
-

The intuitive motivation behind the algorithm is the following. It focuses on reducing the set of candidate hypotheses V , while being careful not to throw away the best classifier $h_{\mathbb{C}}^* = \arg \min_{h \in \mathbb{C}} er(h)$ (supposing, for this informal explanation, that $h_{\mathbb{C}}^*$ exists). Given that this is satisfied at any given time in the algorithm, it makes sense to focus our samples to the region $\text{DIS}(V)$, since a classifier $h_1 \in V$ has smaller error rate than another classifier $h_2 \in V$ if and only if it has smaller conditional error rate given $\text{DIS}(V)$. For this reason, on each round, we seek to remove from V any h for which our confidence bounds indicate that $er(h | \text{DIS}(V)) > er(h_{\mathbb{C}}^* | \text{DIS}(V))$. However, so that we can make use of known results for i.i.d. samples, we freeze the sampling region $R \supseteq \text{DIS}(V)$ and collect an i.i.d. sample from the conditional given this region, updating the region only when doing so allows us to further significantly focus the samples; for this same reason, we also reset the collection of samples Q every time we update the region R , so that it represents samples from the conditional given R . Finally, we maintain the values β_t , which represent confidence upper bounds on $er(h_t) - \nu = (er(h_t | R) - er(h_{\mathbb{C}}^* | R))\mathbb{P}(R)$, and we return the h_t minimizing this confidence bound; note that it does not suffice to return h_n , since the final Q set might be small.

As long as the confidence bounds UB and LB satisfy (overloading notation in the natural way)

$$\mathbb{P}_{Z \sim \mathcal{D}^m} \{\forall h \in \mathbb{C}, LB(h, Z, \delta') \leq er_{\mathcal{D}}(h) \leq UB(h, Z, \delta')\} \geq 1 - \delta'$$

for any distribution \mathcal{D} over $\mathcal{X} \times \mathcal{Y}$ and any $\delta' \in (0, 1)$, and UB and LB converge to each other as m grows, it is known that a $1 - \delta$ confidence bound on $er(\hat{h}_n) - \nu$ converges to 0 [6]. For instance, Balcan, Beygelzimer

and Langford [6] suggest defining these functions based on classic results on uniform convergence rates in passive learning [30], such as

$$(1) \quad \begin{aligned} UB(h, Q, \delta') &= \min\{er_Q(h) + G(|Q|, \delta'), 1\}, \\ LB(h, Q, \delta') &= \max\{er_Q(h) - G(|Q|, \delta'), 0\}, \end{aligned}$$

where $G(m, \delta') = \frac{1}{m} + \sqrt{\frac{\ln(4/\delta') + d \ln(2em/d)}{m}}$ for $m \geq d$, and by convention $G(m, \delta') = \infty$ for $m < d$. This choice of UB and LB is motivated by the following lemma, due to Vapnik [31].

LEMMA 4. *For any distribution \mathcal{D} over $\mathcal{X} \times \mathcal{Y}$, and any $\delta' \in (0, 1)$ and $m \in \mathbb{N}$, with probability $\geq 1 - \delta'$ over the draw of $Z \sim \mathcal{D}^m$, every $h \in \mathbb{C}$ satisfies*

$$(2) \quad |er_Z(h) - er_{\mathcal{D}}(h)| \leq G(m, \delta').$$

To avoid computational issues, instead of explicitly representing the sets V and R , we may implicitly represent them as a set of constraints imposed by the condition in step 7 of previous iterations. We may also replace $\mathbb{P}(\text{DIS}(V))$ and $\mathbb{P}(R)$ by estimates, since these quantities can be estimated to arbitrary precision with arbitrarily high confidence using only *unlabeled* data. Specifically, the convergence rates proven below can be preserved up to constant factors by replacing these quantities with confidence bounds based on a finite number of unlabeled data points; the details of this are included in Appendix C of the supplementary material [20]. As for the number of unlabeled data points required by the above algorithm itself, note that if $\mathbb{P}(\text{DIS}(V))$ becomes small, it will use a large number of unlabeled data points; however, $\mathbb{P}(\text{DIS}(V))$ being small also indicates $er(\hat{h}_n) - \nu$ is small (and indeed β_t). In particular, to get an excess error rate of ε , the algorithm will generally require a number of unlabeled data points only polynomial in $1/\varepsilon$; also, the condition in step 4 guarantees the total number of unlabeled data points used by the algorithm is bounded with high probability. For comparison, recall that passive learning typically requires a number of *labeled* data points polynomial in $1/\varepsilon$.

3.2. Algorithm 2. The second noise-robust algorithm we study was originally proposed by Dasgupta, Hsu and Monteleoni [14]. It uses a type of constrained passive learning subroutine, LEARN , defined as follows for two sets of labeled data points, \mathcal{L} and Q .

$$\text{LEARN}_{\mathbb{C}}(\mathcal{L}, Q) = \arg \min_{h \in \mathbb{C} : er_{\mathcal{L}}(h) = 0} er_Q(h).$$

By convention, if no $h \in \mathbb{C}$ has $er_{\mathcal{L}}(h) = 0$, $\text{LEARN}_{\mathbb{C}}(\mathcal{L}, Q) = \emptyset$. The algorithm is formally defined below, in terms of a sequence of estimators Δ_m , defined later.

Algorithm 2Input: hypothesis class \mathbb{C} , label budget n , confidence δ , functions Δ_m Output: classifier \hat{h}_n , sets of (index, label) pairs \mathcal{L} and Q

-
0. $\mathcal{L} \leftarrow \emptyset, Q \leftarrow \emptyset$
 1. For $m = 1, 2, \dots$
 2. If $|Q| = n$ or $m > 2^n$, Return $\hat{h}_n = \text{LEARN}_{\mathbb{C}}(\mathcal{L}, Q)$ along with \mathcal{L} and Q
 3. For each $y \in \{-1, +1\}$, let $h^{(y)} = \text{LEARN}_{\mathbb{C}}(\mathcal{L} \cup \{(m, y)\}, Q)$
 4. If some y has $h^{(-y)} = \emptyset$ or

$$er_{\mathcal{L} \cup Q}(h^{(-y)}) - er_{\mathcal{L} \cup Q}(h^{(y)}) > \Delta_{m-1}(\mathcal{L}, Q, h^{(y)}, h^{(-y)}, \delta)$$
 5. Then $\mathcal{L} \leftarrow \mathcal{L} \cup \{(m, y)\}$
 6. Else Request the label Y_m and let $Q \leftarrow Q \cup \{(m, Y_m)\}$
-

The algorithm maintains two sets of labeled data points: \mathcal{L} and Q . The set Q represents points of which we have requested the labels. The set \mathcal{L} represents the remaining points, and the labels of points in \mathcal{L} are *inferred*. Specifically, suppose (inductively) that at some time m we have that every $(i, y) \in \mathcal{L}$ has $h_{\mathbb{C}}^*(X_i) = y$, where $h_{\mathbb{C}}^* = \arg \min_{h \in \mathbb{C}} er(h)$ (supposing the min is achieved, for this informal motivation). At any point, we can be fairly confident that $h_{\mathbb{C}}^*$ will have relatively small empirical error rate. Thus, if all of the classifiers h with $er_{\mathcal{L}}(h) = 0$ and $h(X_m) = -y$ have relatively large empirical error rates compared to some h with $er_{\mathcal{L}}(h) = 0$ and $h(X_m) = y$, we can confidently infer that $h_{\mathbb{C}}^*(X_m) = y$. Note that this is not the *true* label Y_m , but a sort of “denoised” version of it. Once we infer this label, since we are already confident that this is the $h_{\mathbb{C}}^*$ label, and $h_{\mathbb{C}}^*$ is the classifier we wish to compete with, we simply add this label as a *constraint*: that is, we require every classifier under consideration in the future to have $h(X_m) = h_{\mathbb{C}}^*(X_m)$. This is how elements of \mathcal{L} are added. On the other hand, if we cannot confidently infer $h_{\mathbb{C}}^*(X_m)$, because some classifiers labeling X_m as $-h_{\mathbb{C}}^*(X_m)$ also have relatively small empirical error rates, then we simply request the label Y_m and add it to the set Q . Note that in order to make this comparison, we needed to be able to calculate the differences of empirical error rates; however, as long as we only consider the set of classifiers h that *agree* on the labels in \mathcal{L} , we will have $er_{\mathcal{L} \cup Q}(h_1) - er_{\mathcal{L} \cup Q}(h_2) = er_m(h_1) - er_m(h_2)$, for any two such classifiers h_1 and h_2 , where $m = |\mathcal{L} \cup Q|$.

The key to the above argument is carefully choosing a threshold for how large the difference in empirical error rates needs to be before we can confidently infer the label. For this purpose, Algorithm 2 is defined in terms of a function, $\Delta_m(\mathcal{L}, Q, h^{(y)}, h^{(-y)}, \delta)$, representing a threshold for a type of hypothesis test. This threshold must be set carefully, since the sequence of labeled data points corresponding to $\mathcal{L} \cup Q$ is not actually an i.i.d. sample from \mathcal{D}_{XY} . Dasgupta, Hsu and Monteleoni [14] suggest defining this function

as

$$(3) \quad \Delta_m(\mathcal{L}, Q, h^{(y)}, h^{(-y)}, \delta) = \beta_m^2 + \beta_m(\sqrt{er_{\mathcal{L} \cup Q}(h^{(y)})} + \sqrt{er_{\mathcal{L} \cup Q}(h^{(-y)})}),$$

where $\beta_m = \sqrt{\frac{4 \ln(8m(m+1)\mathcal{S}(\mathbb{C}, 2m)^2/\delta)}{m}}$ and $\mathcal{S}(\mathbb{C}, 2m)$ is the shatter coefficient (e.g., [15, 31]); this suggestion is based on a confidence bound they derive, and they prove the correctness of the algorithm with this definition, meaning that the $1 - \delta$ confidence bound on its error rate converges to ν as $n \rightarrow \infty$. For now we will focus on the first return value (the classifier), leaving the others for Section 5, where they will be useful for chaining multiple executions together.

4. Convergence rates. In both of the above cases, one can prove guarantees stating that neither algorithm’s convergence rates are ever significantly worse than passive learning by empirical risk minimization [6, 14]. However, it is even more interesting to discuss situations in which one can prove error rate guarantees for these algorithms significantly *better* than those achievable by passive learning. In this section, we begin by reviewing known results on these potential improvements, stated in terms of the disagreement coefficient; we then proceed to discuss new results for Algorithm 1 and a novel variant of Algorithm 2, and describe the convergence rates achieved by these methods in terms of the disagreement coefficient and Tsybakov’s noise conditions.

To simplify the presentation, for the remainder of this paper we will restrict the discussion to situations with $\theta > 0$ (and therefore \mathbb{C} with $d > 0$ too). Handling the extra case of $\theta = 0$ is a trivial matter, since $\theta = 0$ would imply that any proper learning algorithm achieves excess error 0 for all values of n .

4.1. The disagreement coefficient and active learning: Basic results. Before going into the results for general distributions \mathcal{D}_{XY} on $\mathcal{X} \times \mathcal{Y}$, it will be instructive to first look at the special case when the noise rate is zero. Understanding how the disagreement coefficient enters into the analysis of this simpler case may aid in digestion of the theorems and proofs for the general case presented later, where it plays an essentially analogous role. Most of the major ingredients of the proofs for the general case can be found in this special case, albeit in a much simpler form. Although this result has not previously been published, the proof is essentially analogous to (one case of) the analysis of Algorithm 1 in [19].

THEOREM 1. *Let $f \in \mathbb{C}$ be such that $er(f) = 0$ and $\theta_f < \infty$. $\forall n \in \mathbb{N}$ and $\delta \in (0, 1)$, with probability $\geq 1 - \delta$ over the draw of the unlabeled data, the*

classifier \hat{h}_n returned by Algorithm 0 after n label requests satisfies

$$er(\hat{h}_n) \leq 2 \cdot \exp \left\{ -\frac{n}{12\theta_f(d \ln(22\theta_f) + \ln(3n/\delta))} \right\}.$$

PROOF. As in the algorithm, let V_t denote the set of classifiers in \mathbb{C} consistent with the first t label requests. If $\mathbb{P}(\text{DIS}(V_t)) > 0$ for all values of t in the algorithm, then with probability 1 the algorithm uses all n label requests. Technically, each claim below should be followed by the phrase, “unless $\mathbb{P}(\text{DIS}(V_t)) = 0$ for some $t \leq n$, in which case $er(\hat{h}_n) = 0$ so the bound trivially holds.” However, to simplify the presentation, we will make this special case implicit, and will not mention it further.

The high-level outline of this proof is to use $\mathbb{P}(\text{DIS}(V_t))$ as an upper bound on $\sup_{h \in V_t} er(h)$, and then show $\mathbb{P}(\text{DIS}(V_t))$ is halved roughly every $\lambda = \tilde{O}(\theta_f d)$ label requests. Thus, after roughly $\tilde{O}(\theta_f d \log(1/\varepsilon))$ label requests, any $h \in V_t$ should have $er(h) \leq \varepsilon$.

Specifically, let $\lambda_n = \lceil 8\theta_f(d \ln(8e\theta_f) + \ln(2n/\delta)) \rceil$. If $n \leq \lambda_n$, the bound in the theorem statement trivially holds, since the right-hand side exceeds 1; otherwise, consider some nonnegative $t \leq n - \lambda_n$ and $t' = t + \lambda_n$. Let X_{m_t} denote the point corresponding to the t th label request, and let $X_{m_{t'}}$ denote the point corresponding to label request number t' . It must be that

$$|\{X_{m_t+1}, X_{m_t+2}, \dots, X_{m_{t'}}\} \cap \text{DIS}(V_t)| \geq \lambda_n,$$

which means there is an i.i.d. sample of size λ_n , with distribution equivalent to the conditional of X given $\{X \in \text{DIS}(V_t)\}$, contained in $\{X_{m_t+1}, \dots, X_{m_{t'}}\}$: namely, the first λ_n points in this subsequence that are in $\text{DIS}(V_t)$.

Now recall that, by classic results from the passive learning literature (e.g., [5]), this implies that on an event $E_{\delta,t}$ holding with probability $1 - \delta/n$,

$$\sup_{h \in V_{t'}} er(h | \text{DIS}(V_t)) \leq 2 \frac{d \ln(2e\lambda_n/d) + \ln(2n/\delta)}{\lambda_n}.$$

Also note that λ_n was defined (with express purpose) so that

$$2 \frac{d \ln(2e\lambda_n/d) + \ln(2n/\delta)}{\lambda_n} \leq 1/(2\theta_f).$$

Recall that, since $er(f) = 0$, we have $er(h) = \mathbb{P}(h(X) \neq f(X))$. Since $f \in V_{t'} \subseteq V_t$, this means for any $h \in V_{t'}$ we have $\{x : h(x) \neq f(x)\} \subseteq \text{DIS}(V_t)$, and thus

$$\begin{aligned} \sup_{h \in V_{t'}} \mathbb{P}(h(X) \neq f(X)) &= \sup_{h \in V_{t'}} \mathbb{P}(h(X) \neq f(X) | X \in \text{DIS}(V_t)) \mathbb{P}(\text{DIS}(V_t)) \\ &= \sup_{h \in V_{t'}} er(h | \text{DIS}(V_t)) \mathbb{P}(\text{DIS}(V_t)) \leq \mathbb{P}(\text{DIS}(V_t)) / (2\theta_f). \end{aligned}$$

So $V_{t'} \subseteq B(f, \mathbb{P}(\text{DIS}(V_t))/(2\theta_f))$, and therefore by monotonicity of $\mathbb{P}(\text{DIS}(\cdot))$ and the definition of θ_f

$$\mathbb{P}(\text{DIS}(V_{t'})) \leq \mathbb{P}(\text{DIS}(B(f, \mathbb{P}(\text{DIS}(V_t))/(2\theta_f)))) \leq \mathbb{P}(\text{DIS}(V_t))/2.$$

By a union bound, $E_{\delta,t}$ holds for every $t \in \{i\lambda_n : i \in \{0, 1, \dots, \lfloor n/\lambda_n \rfloor - 1\}\}$ with probability $\geq 1 - \delta$. On these events, if $n \geq \lambda_n \lceil \log_2(1/\varepsilon) \rceil$, then (by induction)

$$\sup_{h \in V_n} \text{er}(h) \leq \mathbb{P}(\text{DIS}(V_n)) \leq \varepsilon.$$

Solving for ε in terms of n gives the result (with a slight increase in constants due to relaxing the ceiling functions). \square

4.2. Known results on convergence rates for agnostic active learning. We will now describe the known results for agnostic active learning algorithms, starting with Algorithm 1. The key to the potential convergence rate improvements of Algorithm 1 is that, as the region of disagreement R decreases in measure, the error difference $\text{er}(h|R) - \text{er}(h'|R)$ of any classifiers $h, h' \in V$ under the *conditional* sampling distribution (given R) can become significantly larger [by a factor of $\mathbb{P}(R)^{-1}$] than $\text{er}(h) - \text{er}(h')$, making it significantly easier to determine which of the two is worse using a sample of labeled data. In particular, [19] developed a technique for analyzing this type of algorithm, and adapting that analysis to the above definition of Algorithm 1 results in the following guarantee.

THEOREM 2 [19]. *Let \hat{h}_n be the classifier returned by Algorithm 1 when allowed n label requests, using the bounds (1) and confidence parameter $\delta \in (0, 1/2)$. Then there exists a finite universal constant c such that, with probability $\geq 1 - \delta$, $\forall n \in \mathbb{N}$,*

$$\begin{aligned} \text{er}(\hat{h}_n) - \nu &\leq c \sqrt{\frac{\nu^2 \theta^2 (d \log n + \log(1/\delta)) \log((n + 2\nu\theta)/(\nu\theta))}{n}} \\ &\quad + 2 \exp \left\{ - \frac{n}{c \theta^2 (d \log \theta + \log(n/\delta))} \right\}. \end{aligned}$$

Similarly, the key to improvements from Algorithm 2 is that as the number m of processed unlabeled data points increases, we only need to request the labels of those data points in the region of disagreement of the set of classifiers with near-optimal empirical error rates. Thus, if the region of disagreement of classifiers with excess error $\leq \varepsilon$ shrinks as ε shrinks, we expect the frequency of label requests to shrink as m increases. Since we are careful not to discard the best classifier, and the excess error rate of a classifier can be bounded in terms of the Δ_m function, we end up

with a bound on the excess error which is converging in m , the number of *unlabeled* data points processed, even though we request a number of labels growing slower than m . When this situation occurs, we expect Algorithm 2 will provide an improved convergence rate compared to passive learning. Dasgupta, Hsu and Monteleoni [14] prove the following convergence rate guarantee.

THEOREM 3 [14]. *Let \hat{h}_n be the classifier returned by Algorithm 2 when allowed n label requests, using the threshold (3), and confidence parameter $\delta \in (0, 1/2)$. Then there exists a finite universal constant c such that, with probability $\geq 1 - \delta$, $\forall n \in \mathbb{N}$,*

$$\begin{aligned} \text{er}(\hat{h}_n) - \nu &\leq c \sqrt{\frac{\nu^2 \theta (d \log((n + 2\nu\theta)/(\nu\theta)) + \log(1/\delta))}{n}} \\ &\quad + c \left(d + \log \frac{1}{\delta} \right) \exp \left\{ - \sqrt{\frac{n}{c\theta(d + \log(1/\delta))}} \right\}. \end{aligned}$$

Note that, among other changes, this bound improves the dependence on the disagreement coefficient θ , compared to the bound for Algorithm 1. In both cases, for certain ranges of θ , ν and n , these bounds can represent significant improvements in the excess error guarantees, compared to the corresponding guarantees possible for passive learning. However, in both cases, when $\nu > 0$ these bounds have an *asymptotic* dependence on n of $\tilde{\Theta}(n^{-1/2})$, which is no better than the convergence rates achievable by passive learning (e.g., by empirical risk minimization). Thus, there remains the question of whether either algorithm can achieve asymptotic convergence rates strictly superior to passive learning for distributions with nonzero noise rates. This is the topic we turn to next.

4.3. Active learning under Tsybakov's noise conditions. It is known that for most nontrivial \mathbb{C} , for any n and $\nu > 0$, for every active learning algorithm there is some distribution with noise rate ν for which we can guarantee excess error no better than $\propto \nu n^{-1/2}$ [21]; that is, the $n^{-1/2}$ asymptotic dependence on n in the above bounds matches the corresponding minimax rate, and thus cannot be improved as long as the bounds depend on \mathcal{D}_{XY} only via ν (and θ). Therefore, if we hope to discover situations in which these algorithms have strictly superior asymptotic dependence on n , we will need to allow the bounds to depend on a more detailed description of the noise distribution than simply the noise rate ν .

As previously mentioned, one way to describe a noise distribution using a more detailed parametrization is to use Tsybakov's noise conditions (Condition 1). In the context of passive learning, this allows one to describe

situations in which the rate of convergence is between n^{-1} and $n^{-1/2}$, even when $\nu > 0$. This raises the natural question of how these active learning algorithms perform when the noise distribution satisfies this condition with finite μ and κ parameter values. In many ways, it seems active learning is particularly well-suited to exploit these more favorable noise conditions, since they imply that as we eliminate suboptimal classifiers, the diameter of the remaining set shrinks; thus, for finite θ values, the region of disagreement should also be shrinking, allowing us to focus the samples in a smaller region and accelerate the convergence.

Focusing on the special case of learning one-dimensional threshold classifiers under a certain uniform marginal distribution, Castro and Nowak [12] studied conditions related to Condition 1. In particular, they studied a threshold-learning algorithm that, unlike the algorithms described here, takes κ as *input*, and found its convergence rate to be $\propto (\frac{\log n}{n})^{\kappa/(2\kappa-2)}$ when $\kappa > 1$, and $\exp\{-cn\}$ for some (μ -dependent) constant c , when $\kappa = 1$. Note that this improves over the $n^{-\kappa/(2\kappa-1)}$ rates achievable in passive learning [12, 28]. Subsequently, Balcan, Broder and Zhang [7] proved an analogous positive result for higher-dimensional linear separators, and Castro and Nowak [12] additionally showed a related result for boundary fragment classes (see below); in both cases, the algorithm depends explicitly on the noise parameters. Later, in a personal communication, Langford and Castro claimed that in fact Algorithm 1 achieves this rate (up to log factors) for the one-dimensional thresholds problem, leading to speculation that perhaps these improvements are achievable in the general case as well (under conditions on the disagreement coefficient). Castro and Nowak [12] also prove that a value $\propto n^{-\kappa/(2\kappa-2)}$ (or $\exp\{-c'n\}$, for some c' , when $\kappa = 1$) is also a *lower bound* on the minimax rate for the threshold learning problem. In fact, a similar proof to theirs can be used to show this same lower bound holds for any nontrivial \mathbb{C} . For completeness, a proof of this more general result is included in Appendix D of the supplementary material [20].

Other than the few specific results mentioned above, it was not previously known whether Algorithm 1 or Algorithm 2, or indeed *any* active learning algorithm, generally achieves convergence rates that exhibit these types of improvements.

4.4. *Adaptive rates in active learning: Algorithm 1.* The above observations open the question of whether these algorithms, or variants thereof, improve this asymptotic dependence on n . It turns out this is indeed possible. Specifically, we have the following result for Algorithm 1.

THEOREM 4. *Let \hat{h}_n be the classifier returned by Algorithm 1 when allowed n label requests, using the bounds (1) and confidence parameter $\delta \in (0, 1/2)$. Suppose further that \mathcal{D}_{XY} satisfies Condition 1. Then there*

exists a finite (κ - and μ -dependent) constant c such that, for any $n \in \mathbb{N}$, with probability $\geq 1 - \delta$,

$$er(\hat{h}_n) - \nu \leq \begin{cases} 2 \cdot \exp\left\{-\frac{n}{c\theta^2(d \log n + \log(1/\delta))}\right\}, & \text{when } \kappa = 1, \\ c \left(\frac{\theta^2(d \log n + \log(1/\delta)) \log n}{n}\right)^{\kappa/(2\kappa-2)}, & \text{when } \kappa > 1. \end{cases}$$

PROOF. We will proceed by bounding the *label complexity*, or size of the label budget n that is sufficient to guarantee, with high probability, that the excess error of the returned classifier will be at most ε (for arbitrary $\varepsilon > 0$); with this in hand, we can simply bound the inverse of the function to get the result in terms of a bound on excess error.

Throughout this proof (and proofs of later results in this paper), we will make frequent use of basic facts about $er(h|R)$. In particular, for any classifiers h, h' and set $R \subseteq \mathcal{X}$, we have $er(h) = er(h|R)\mathbb{P}(R) + er(h|\mathcal{X} \setminus R)\mathbb{P}(\mathcal{X} \setminus R)$; also, if $\{x : h(x) \neq h'(x)\} \subseteq R$, we have $er(h|\mathcal{X} \setminus R) - er(h'|\mathcal{X} \setminus R) = 0$ and therefore $er(h) - er(h') = (er(h|R) - er(h'|R))\mathbb{P}(R)$.

Note that, by Lemma 4 and a union bound, on an event of probability $1 - \delta$, (2) holds with $\delta' = \delta/n$ for every set Q , relative to the conditional distribution given its respective R set, for any value of n . For the remainder of this proof, we assume that this $1 - \delta$ probability event occurs. In particular, this means that for every $h \in \mathbb{C}$ and every Q set in the algorithm, $LB(h, Q, \delta/n) \leq er(h|R) \leq UB(h, Q, \delta/n)$, for the set R that Q is sampled under.

Our first task is to show that we never remove the “good” classifiers from V . We only remove a classifier h from V if $h' = \arg \min_{h' \in V} UB(h', Q, \delta/n)$ has $LB(h, Q, \delta/n) > UB(h', Q, \delta/n)$. Each $h \in V$ has $\{x : h(x) \neq h'(x)\} \subseteq \text{DIS}(V) \subseteq R$, so that

$$UB(h', Q, \delta/n) - LB(h, Q, \delta/n) \geq er(h'|R) - er(h|R) = \frac{er(h') - er(h)}{\mathbb{P}(R)}.$$

Thus, for any $h \in V$ with $er(h) \leq er(h')$, $UB(h', Q, \delta/n) - LB(h, Q, \delta/n) \geq er(h'|R) - er(h|R) = (er(h') - er(h))/\mathbb{P}(R) \geq 0$, so that on any given round of the algorithm, the set $\{h \in V : er(h) \leq er(h')\}$ is not removed from V . In particular, since we always have $er(h') \geq \nu$, by induction this implies the invariant $\inf_{h \in V} er(h) = \nu$, and therefore also

$$\begin{aligned} \forall t \quad er(h_t) - \nu &= er(h_t) - \inf_{h \in V} er(h) \\ &= \left(er(h_t|R) - \inf_{h \in V} er(h|R) \right) \mathbb{P}(R) \leq \beta_t, \end{aligned}$$

where again the second equality is due to the fact that $\forall h \in V, \{x : h_t(x) \neq h(x)\} \subseteq \text{DIS}(V) \subseteq R$. We will spend the remainder of the proof bounding the size of n sufficient to guarantee some $\beta_t \leq \varepsilon$. In particular, similar to the proof of Theorem 1, we will see that as long as $\beta_t > \varepsilon$, we will halve $\mathbb{P}(\text{DIS}(V))$ roughly every $\tilde{O}(\theta^2 d \varepsilon^{2/\kappa-2})$ label requests, so that the total number of label requests before some $\beta_t \leq \varepsilon$ is at most roughly $\tilde{O}(\theta^2 d \varepsilon^{2/\kappa-2} \log(1/\varepsilon))$.

Recalling the definition of $h^{[k]}$ (from Definition 2), let

$$(4) \quad V^{(\theta)} = \left\{ h \in V : \limsup_{k \rightarrow \infty} \mathbb{P}(h(X) \neq h^{[k]}(X)) > \frac{\mathbb{P}(R)}{2\theta} \right\}.$$

Note that after step 7, if $V^{(\theta)} = \emptyset$, then

$$\begin{aligned} \mathbb{P}(\text{DIS}(V)) &\leq \mathbb{P}\left(\text{DIS}\left(\left\{h \in \mathbb{C} : \limsup_{k \rightarrow \infty} \mathbb{P}(h(X) \neq h^{[k]}(X)) \leq \mathbb{P}(R)/(2\theta)\right\}\right)\right) \\ &= \lim_{k' \rightarrow \infty} \mathbb{P}\left(\text{DIS}\left(\bigcap_{k > k'} B(h^{[k]}, \mathbb{P}(R)/(2\theta))\right)\right) \\ &\leq \lim_{k' \rightarrow \infty} \mathbb{P}\left(\bigcap_{k > k'} \text{DIS}(B(h^{[k]}, \mathbb{P}(R)/(2\theta)))\right) \\ &\leq \liminf_{k \rightarrow \infty} \mathbb{P}(\text{DIS}(B(h^{[k]}, \mathbb{P}(R)/(2\theta)))) \\ &\leq \liminf_{k \rightarrow \infty} \theta_{h^{[k]}} \frac{\mathbb{P}(R)}{2\theta} = \frac{\mathbb{P}(R)}{2}, \end{aligned}$$

so that we will satisfy the condition in step 2 on the next round. Here we have used the definition of θ in the final inequality and equality. On the other hand, if after step 7, we have $V^{(\theta)} \neq \emptyset$, then

$$\begin{aligned} \emptyset &\neq \left\{ h \in V : \limsup_{k \rightarrow \infty} \mathbb{P}(h(X) \neq h^{[k]}(X)) > \frac{\mathbb{P}(R)}{2\theta} \right\} \\ &= \left\{ h \in V : \left(\frac{\limsup_{k \rightarrow \infty} \mathbb{P}(h(X) \neq h^{[k]}(X))}{\mu} \right)^\kappa > \left(\frac{\mathbb{P}(R)}{2\mu\theta} \right)^\kappa \right\} \\ &\subseteq \left\{ h \in V : \left(\frac{\text{diam}(er(h) - \nu; \mathbb{C})}{\mu} \right)^\kappa > \left(\frac{\mathbb{P}(R)}{2\mu\theta} \right)^\kappa \right\} \\ &\subseteq \left\{ h \in V : er(h) - \nu > \left(\frac{\mathbb{P}(R)}{2\mu\theta} \right)^\kappa \right\} \\ &= \left\{ h \in V : er(h|R) - \inf_{h' \in V} er(h'|R) > \mathbb{P}(R)^{\kappa-1} (2\mu\theta)^{-\kappa} \right\} \\ &\subseteq \left\{ h \in V : UB(h, Q, \delta/n) - \min_{h' \in V} LB(h', Q, \delta/n) > \mathbb{P}(R)^{\kappa-1} (2\mu\theta)^{-\kappa} \right\} \end{aligned}$$

$$\begin{aligned} &\subseteq \left\{ h \in V : LB(h, Q, \delta/n) - \min_{h' \in V} UB(h', Q, \delta/n) \right. \\ &\quad \left. > \mathbb{P}(R)^{\kappa-1} (2\mu\theta)^{-\kappa} - 4G(|Q|, \delta/n) \right\}. \end{aligned}$$

Here, the third line follows from the fact that $er(h^{[k]}) \leq er(h)$ for all sufficiently large k , the fourth line follows from Condition 1, and the final line follows from the definition of UB and LB . By definition, every $h \in V$ has $LB(h, Q, \delta/n) \leq \min_{h' \in V} UB(h', Q, \delta/n)$, so for this last set to be nonempty after step 7, we must have $\mathbb{P}(R)^{\kappa-1} (2\mu\theta)^{-\kappa} < 4G(|Q|, \delta/n)$.

Combining these two cases ($V^{(\theta)} = \emptyset$ and $V^{(\theta)} \neq \emptyset$), since $|Q|$ gets reset to 0 upon reaching step 3, we have that after every execution of step 7,

$$(5) \quad \mathbb{P}(R)^{\kappa-1} (2\mu\theta)^{-\kappa} < 4G(|Q| - 1, \delta/n).$$

If $\mathbb{P}(R) \leq \frac{\varepsilon}{2G(|Q|-1, \delta/n)} \leq \frac{\varepsilon}{2G(|Q|, \delta/n)}$, then certainly $\beta_t \leq \varepsilon$ (by definition of $\beta_t \leq 2G(|Q|, \delta/n)\mathbb{P}(R)$). So on any round for which $\beta_t > \varepsilon$, we must have

$$(6) \quad \frac{\varepsilon}{2G(|Q| - 1, \delta/n)} < \mathbb{P}(R).$$

Combining (5) and (6), on any round for which $\beta_t > \varepsilon$,

$$(7) \quad \left(\frac{\varepsilon}{2G(|Q| - 1, \delta/n)} \right)^{\kappa-1} (2\mu\theta)^{-\kappa} < 4G(|Q| - 1, \delta/n).$$

Solving for $G(|Q| - 1, \delta/n)$ reveals that when $\beta_t > \varepsilon$,

$$(8) \quad 4^{-1/\kappa} \left(\frac{\varepsilon}{2} \right)^{(\kappa-1)/\kappa} (2\mu\theta)^{-1} < G(|Q| - 1, \delta/n).$$

Basic algebra shows that when $n \geq |Q| > d$, we have

$$G(|Q| - 1, \delta/n) \leq 3\sqrt{\frac{\ln(4/\delta) + (d+1)\ln(n)}{|Q|}}.$$

Combining this with (8), solving for $|Q|$ and adding d to handle the case $|Q| \leq d$, we have that on any round for which $\beta_t > \varepsilon$,

$$(9) \quad |Q| \leq \left(\frac{2}{\varepsilon} \right)^{(2\kappa-2)/\kappa} (6\mu\theta)^2 4^{2/\kappa} \left(\ln \frac{4}{\delta} + (d+1)\ln(n) \right) + d.$$

Since $\beta_t \leq \mathbb{P}(R)$ by definition, and $\mathbb{P}(R)$ is at least halved each time we reach step 3, we need to reach step 3 at most $\lceil \log_2(1/\varepsilon) \rceil$ times before we are guaranteed some $\beta_t \leq \varepsilon$. Thus, any

$$(10) \quad n \geq 1 + \left(\left(\frac{2}{\varepsilon} \right)^{(2\kappa-2)/\kappa} (6\mu\theta)^2 4^{2/\kappa} \left(\ln \frac{4}{\delta} + (d+1)\ln(n) \right) + d \right) \log_2 \frac{2}{\varepsilon}$$

suffices to guarantee either some $|Q|$ exceeds (9) or we reach step 3 at least $\lceil \log_2(1/\varepsilon) \rceil$ times, either of which implies the existence of some $\beta_t \leq \varepsilon$. The stated result now follows by basic inequalities to bound the smallest value of ε satisfying (10) for a given value of n . \square

If the disagreement coefficient is finite, Theorem 4 can often represent a significant improvement in convergence rate compared to passive learning, where we typically expect rates of order $n^{-\kappa/(2\kappa-1)}$ [12, 26, 28]; this gap is especially notable when the disagreement coefficient and κ are small. Furthermore, the bound matches (up to logarithmic factors) the form of the minimax rate *lower bound* proved by Castro and Nowak [12] for threshold classifiers (where $\theta = 2$); as mentioned, that lower bound proof can be generalized to any nontrivial \mathbb{C} (see Appendix D of the supplementary material [20]), so that the rate of Theorem 4 is nearly minimax optimal for any nontrivial \mathbb{C} with *bounded* disagreement coefficients. Also note that, unlike the upper bound analysis of Castro and Nowak [12], we do not require the algorithm to be given any extra information about the noise distribution, so that this result is somewhat stronger; it is also more general, as this bound applies to an arbitrary hypothesis class.

A refined analysis and minor tweaks to the algorithm should be able to reduce the log factors in this result. For instance, defining UB and LB using the uniform convergence bounds of Alexander [1], and using a slightly more complicated algorithm closer to the original definition [6, 19]—taking multiple samples between bound evaluations, allowing a larger confidence argument to the UB and LB evaluations—the $\log^2 n$ factor should reduce at least to $\log n \log \log n$, if not further. Also, as previously mentioned, it is possible to replace the quantities $\mathbb{P}(R)$ and $\mathbb{P}(\text{DIS}(V))$ in Algorithm 1 with estimators of these quantities based on a finite sample of unlabeled data points, while preserving the results of Theorem 4 up to constant factors. We include an example of such estimators in Appendix C of the supplementary material [20], along with a sketch of how to modify the proof of Theorem 4 to compensate for using these estimated probabilities.

4.5. *Adaptive rates in active learning: Algorithm 2.* Note that, as before, n gets divided by θ^2 in the rates achieved by Algorithm 1. As before, it is not clear whether any modification to the definitions of UB and LB can reduce this exponent on θ from 2 to 1. As such, it is natural to investigate the rates achieved by Algorithm 2 under Condition 1; we know that it does improve the dependence on θ for the worst case rates over distributions with any given noise rate, so we might hope that it does the same for the rates over distributions with any given values of μ and κ . Unfortunately, we do not presently know whether the original definition of Algorithm 2 achieves

this improvement. However, we now present a slight modification of the algorithm, and prove that it does indeed provide the desired improvement in dependence on θ , while maintaining the improvements in the asymptotic dependence on n . Specifically, consider the following definition for the threshold in Algorithm 2:

$$(11) \quad \Delta_m(\mathcal{L}, Q, h^{(y)}, h^{(-y)}, \delta) = 3\hat{\mathcal{E}}_{\mathbb{C}}(\mathcal{L} \cup Q, \delta; \mathcal{L}),$$

where $\hat{\mathcal{E}}_{\mathbb{C}}(\cdot, \cdot; \cdot)$ is defined in the [Appendix](#), based on a notion of local Rademacher complexity studied by Koltchinskii [23]. In particular, the quantity $\hat{\mathcal{E}}_{\mathbb{C}}$ is known to be adaptive to Tsybakov's noise conditions, in the sense that more favorable noise conditions yield smaller values of $\hat{\mathcal{E}}_{\mathbb{C}}$. Using this definition, we have the following theorem; due to space limitations, its proof is not presented here, but is included in Appendix B of the supplementary material [20].

THEOREM 5. *Suppose \hat{h}_n is the classifier returned by Algorithm 2 with threshold as in (11), when allowed n label requests and given confidence parameter $\delta \in (0, 1/2)$. Suppose further that \mathcal{D}_{XY} satisfies Condition 1 with finite parameter values κ and μ . Then there exists a finite (κ and μ -dependent) constant c such that, with probability $\geq 1 - \delta$, $\forall n \in \mathbb{N}$,*

$$er(\hat{h}_n) - \nu \leq \begin{cases} c \left(d + \log \frac{1}{\delta} \right) \cdot \exp \left\{ - \sqrt{\frac{n}{c\theta(d + \log(1/\delta))}} \right\}, & \text{when } \kappa = 1, \\ c \left(\frac{\theta(d \log n + \log(1/\delta))}{n} \right)^{\kappa/(2\kappa-2)}, & \text{when } \kappa > 1. \end{cases}$$

Note that this does indeed improve the dependence on θ , reducing its exponent from 2 to 1; we do lose some in that there is now a square root in the exponent of the $\kappa = 1$ case; however, as with Theorem 4, it is likely that slight refinements to the definition of Δ_m would reduce this (though we may also need to weaken the theorem statement to hold for any single n , rather than simultaneously for all n).

The bound in Theorem 5 is stated in terms of the VC dimension d . However, for certain nonparametric hypothesis classes, it is sometimes preferable to quantify the complexity of the class in terms of a constraint on the *entropy* of the class, relative to the distribution \mathcal{D}_{XY} (see e.g., [12, 23, 28, 29]). Specifically, for $\varepsilon \in [0, 1]$, define

$$\omega_{\mathbb{C}}(m, \varepsilon) = \mathbb{E} \sup_{\substack{h_1, h_2 \in \mathbb{C}: \\ \mathbb{P}\{h_1(X) \neq h_2(X)\} \leq \varepsilon}} |(er(h_1) - er_m(h_1)) - (er(h_2) - er_m(h_2))|.$$

CONDITION 2. There exist finite constants $\alpha > 0$ and $\rho \in (0, 1)$ s.t. $\forall m \in \mathbb{N}$ and $\varepsilon \in [0, 1]$, $\omega_{\mathbb{C}}(m, \varepsilon) \leq \alpha \cdot \max\{\varepsilon^{(1-\rho)/2} m^{-1/2}, m^{-1/(1+\rho)}\}$.

In particular, the entropy with bracketing condition used in the original minimax analysis of Tsybakov [28] implies Condition 2 [23], as does the analogous condition for random entropy [17, 18, 24]. In passive learning, it is known that empirical risk minimization achieves a rate of order $n^{-\kappa/(2\kappa+\rho-1)}$ under Conditions 1 and 2 [23, 24] (see also Appendix B of the supplementary material [20], especially (19) and Lemma 5), and that this is sometimes minimax optimal [28]. The following theorem gives a bound on the rate of convergence of the same version of Algorithm 2 as in Theorem 5, this time in terms of the entropy condition which, as before, is faster than the passive learning rate when the disagreement coefficient is finite. The proof of this result is included in Appendix B of the supplementary material [20].

THEOREM 6. Suppose \hat{h}_n is the classifier returned by Algorithm 2 with threshold as in (11), when allowed n label requests and given confidence parameter $\delta \in (0, 1/2)$. Suppose further that \mathcal{D}_{XY} satisfies Condition 1 with finite parameter values κ and μ , and Condition 2 with parameter values α and ρ . Then there exists a finite (κ , μ , α and ρ -dependent) constant c such that, with probability $\geq 1 - \delta$, $\forall n \in \mathbb{N}$,

$$er(\hat{h}_n) - \nu \leq c \left(\frac{\theta \log(n/\delta)}{n} \right)^{\kappa/(2\kappa+\rho-2)}.$$

Again, it is likely that refinements to the Δ_m definition may lead to improvements in the log factor. Also, although this result is stated for Algorithm 2, it is conceivable that, by modifying Algorithm 1 to use definitions of V and β_t based on $\hat{\mathcal{E}}_{\mathbb{C}}(Q, \delta; \emptyset)$, an analogous result might be possible for Algorithm 1 as well.

It is worth mentioning that Castro and Nowak [12] proved a minimax lower bound for the hypothesis class of *boundary fragments*, with an exponent having a similar dependence on related definitions of κ and ρ parameters to that of Theorem 6. Their result does provide a valid lower bound here; however, it is not clear whether their lower bound, Theorem 6, both, or neither is tight in the present context, since the value of θ is not presently known for that particular problem, and the matching upper bound of [12] was proven under a stronger restriction on the noise than Condition 1. However, see [33] for an analysis of the disagreement coefficient for other non-parametric hypothesis classes, characterized by smoothness of the decision surface.

5. Model selection. While the previous sections address adaptation to the noise distribution, they are still restrictive in that they deal with hypothesis classes of limited expressiveness. That is, the assumption of finite VC dimension implies a strong restriction on the variety of classifiers one can represent (or approximate) in the class; the entropy conditions allow slightly more flexibility, but under nontrivial distributions, even the entropy conditions imply a significant restriction on the expressiveness of the class. Thus, for algorithms restricted to classifiers from such a restricted hypothesis class, it is often unrealistic to expect convergence to the Bayes error rate. We address this issue in this section by developing a general algorithm for learning with a sequence of nested hypothesis classes of increasing complexity, similar to the setting of Structural Risk Minimization in passive learning [30]. The objective is to adapt, not only to the noise conditions, but also to the complexity of the optimal classifier. The starting point for this discussion is the assumption of a structure on \mathbb{C} , in the form of a sequence of nested hypothesis classes:

$$\mathbb{C}_1 \subset \mathbb{C}_2 \subset \dots$$

Each class has an associated noise rate $\nu_i = \inf_{h \in \mathbb{C}_i} er(h)$, and we define $\nu_\infty = \lim_{i \rightarrow \infty} \nu_i$. We also let θ_i and d_i be the disagreement coefficient and VC dimension, respectively, for the set \mathbb{C}_i . We are interested in an algorithm that guarantees convergence in probability of the error rate to ν_∞ . We are particularly interested in situations where $\nu_\infty = \nu^*$, a condition which is realistic in this setting since the sets \mathbb{C}_i can be defined so that it is always satisfied, even while maintaining each $d_i < \infty$ (see, e.g., [15]). Additionally, if we are so lucky as to have some $\nu_i = \nu^*$, then we would like the convergence rate achieved by the algorithm to be not significantly worse than running one of the above agnostic active learning algorithms with hypothesis class \mathbb{C}_i alone. In this context, we can define a structure-dependent version of Tsybakov's noise condition as follows.

CONDITION 3. For some nonempty $I \subseteq \mathbb{N}$, for each $i \in I$, there exist finite constants $\mu_i > 0$ and $\kappa_i \geq 1$, such that $\forall \varepsilon > 0, \text{diam}(\varepsilon; \mathbb{C}_i) \leq \mu_i \varepsilon^{1/\kappa_i}$.

Note that we do not require every \mathbb{C}_i , $i \in \mathbb{N}$, to have finite μ_i and κ_i , only some nonempty set $I \subseteq \mathbb{N}$; this is important, since we might not expect \mathbb{C}_i to satisfy Condition 1 for small indices i , where the expressiveness is quite restricted.

In passive learning, there are several methods for this type of model selection which are known to preserve the convergence rates of each class \mathbb{C}_i under Condition 3 (e.g., [23, 28]). In particular, Koltchinskii [23] develops a method that performs this type of model selection; it turns out we can modify Koltchinskii's method to suit our present needs in the context of active

learning; this results in a general active learning model selection method that preserves the types of improved rates discussed in the previous section. This modification is presented below, based on using Algorithm 2 as a subroutine. (It should also be possible to define an analogous method that uses Algorithm 1 as a subroutine instead.)

Algorithm 3

Input: nested sequence of classes $\{\mathbb{C}_i\}$, label budget n , confidence parameter δ

Output: classifier \hat{h}_n

0. For $i = \lfloor \sqrt{n/2} \rfloor, \lfloor \sqrt{n/2} \rfloor - 1, \lfloor \sqrt{n/2} \rfloor - 2, \dots, 1$
 1. Let \mathcal{L}_{in} and Q_{in} be the sets returned by Algorithm 2 run with \mathbb{C}_i and the threshold (11), allowing $\lfloor n/(2i^2) \rfloor$ label requests, and confidence $\delta/(2i^2)$
 2. Let $h_{in} \leftarrow \text{LEARN}_{\mathbb{C}_i}(\bigcup_{j \geq i} \mathcal{L}_{jn}, Q_{in})$
 3. If $h_{in} \neq \emptyset$ and $\forall j$ s.t. $i < j \leq \lfloor \sqrt{n/2} \rfloor$,

$$er_{\mathcal{L}_{jn} \cup Q_{jn}}(h_{in}) - er_{\mathcal{L}_{jn} \cup Q_{jn}}(h_{jn}) \leq \frac{3}{2} \hat{\mathcal{E}}_{\mathbb{C}_j}(\mathcal{L}_{jn} \cup Q_{jn}, \delta/(2j^2); \mathcal{L}_{jn})$$
 4. $\hat{h}_n \leftarrow h_{in}$
5. Return \hat{h}_n

The function $\hat{\mathcal{E}}(\cdot, \cdot; \cdot)$ is defined in the [Appendix](#). This method can be shown to have a confidence bound on its error rate converging to ν_∞ at a rate never significantly worse than the original passive learning method of Koltchinskii [23], as desired. Additionally, we have the following guarantee on the rate of convergence under Condition 3. The proof is similar in style to Koltchinskii's original proof, though some care is needed due to the altered sampling distribution and the constraint set \mathcal{L}_{jn} . The proof is included in Appendix B of the supplementary material [20].

THEOREM 7. *Suppose \hat{h}_n is the classifier returned by Algorithm 3, when allowed n label requests and confidence parameter $\delta \in (0, 1/2)$. Suppose further that \mathcal{D}_{XY} satisfies Condition 3. Then there exist finite (κ_i and μ_i -dependent) constants c_i such that, with probability $\geq 1 - \delta$, $\forall n \in \mathbb{N}$,*

$$\begin{aligned}
 & er(\hat{h}_n) - \nu_\infty \\
 & \leq 3 \min_{i \in I} (\nu_i - \nu_\infty) \\
 & \quad + \begin{cases} c_i \left(d_i + \log \frac{1}{\delta} \right) \cdot \exp \left\{ - \sqrt{\frac{n}{c_i \theta_i (d_i + \log(1/\delta))}} \right\}, & \text{if } \kappa_i = 1, \\ c_i \left(\frac{\theta_i (d_i \log n + \log(1/\delta))}{n} \right)^{\kappa_i / (2\kappa_i - 2)}, & \text{if } \kappa_i > 1. \end{cases}
 \end{aligned}$$

In particular, if we are so lucky as to have $\nu_i = \nu^*$ for some finite i , then the above algorithm achieves a convergence rate not significantly worse than that guaranteed by Theorem 5 for applying Algorithm 2 directly, with hypothesis class \mathbb{C}_i . Note that the algorithm itself has no dependence on the set I , nor has it any dependence on each class's complexity parameters $d_i, \kappa_i, \mu_i, \theta_i$; the adaptive behavior of the data-dependent bound $\hat{\mathcal{E}}_{\mathbb{C}_j}$ allows the algorithm to adaptively ignore the returned classifier from the runs of Algorithm 2 for which convergence is slow, thus automatically selecting an index for which the error rate is relatively small.

As in the previous section, we can also show a variant of this result when the complexities are quantified in terms of the entropy. Specifically, consider the following condition and theorem; the proof is in Appendix B of the supplementary material [20]. Again, this represents an improvement over known results for passive learning when the disagreement coefficients are finite.

CONDITION 4. For each $i \in \mathbb{N}$, there exist finite constants $\alpha_i > 0$, $\rho_i \in (0, 1)$ s.t. $\forall m \in \mathbb{N}$ and $\varepsilon \in [0, 1]$, $\omega_{\mathbb{C}_i}(m, \varepsilon) \leq \alpha_i \cdot \max\{\varepsilon^{(1-\rho_i)/2} m^{-1/2}, m^{-1/(1+\rho_i)}\}$.

THEOREM 8. Suppose \hat{h}_n is the classifier returned by Algorithm 3, when allowed n label requests and confidence parameter $\delta \in (0, 1/2)$. Suppose further that \mathcal{D}_{XY} satisfies Conditions 3 and 4. Then there exist finite ($\kappa_i, \mu_i, \alpha_i$ and ρ_i -dependent) constants c_i such that, with probability $\geq 1 - \delta$, $\forall n \in \mathbb{N}$,

$$er(\hat{h}_n) - \nu_\infty \leq 3 \min_{i \in I} (\nu_i - \nu_\infty) + c_i \left(\frac{\theta_i \log(n/\delta)}{n} \right)^{\kappa_i / (2\kappa_i + \rho_i - 2)}.$$

In addition to these theorems for this structure-dependent version of Tsybakov's noise conditions, we also have the following result for a structure-independent noise condition, in the sense that the noise condition does not depend on the particular choice of \mathbb{C}_i sets, but only on the distribution \mathcal{D}_{XY} (and in some sense, the full class $\mathbb{C} = \bigcup_i \mathbb{C}_i$); it may be particularly useful when the class \mathbb{C} is universal, in the sense that it can approximate any classifier.

THEOREM 9. Suppose the sequence $\{\mathbb{C}_i\}$ is constructed so that $\nu_\infty = \nu^*$, and \hat{h}_n is the classifier returned by Algorithm 3, when allowed n label requests and confidence parameter $\delta \in (0, 1/2)$. Suppose that there exists a constant $\mu > 0$ s.t. for all measurable $h: \mathcal{X} \rightarrow \mathcal{Y}$, $er(h) - \nu^* \geq \mu \mathbb{P}\{h(X) \neq h^*(X)\}$. Then there exists a finite (μ -dependent) constant c such that, with probability $\geq 1 - \delta$, $\forall n \in \mathbb{N}$,

$$er(\hat{h}_n) - \nu^* \leq c \min_{i \in \mathbb{N}} (\nu_i - \nu^*) + \left(d_i + \log \frac{i}{\delta} \right) \cdot \exp \left\{ - \sqrt{\frac{n}{c i^2 \theta_i (d_i + \log(i/\delta))}} \right\}.$$

The condition $\nu_\infty = \nu^*$ is quite easy to satisfy: for example, \mathbb{C}_i could be axis-aligned decision trees of depth i , or thresholded polynomials of degree i , or multi-layer neural networks with i internal units, etc. As for the noise condition in Theorem 9, this would be satisfied whenever $\mathbb{P}(|\eta(X) - 1/2| \geq c) = 1$ for some constant $c \in (0, 1/2]$. The case where $er(h) - \nu^* \geq \mu\mathbb{P}\{h(X) \neq h^*(X)\}^\kappa$ for $\kappa > 1$ can be studied analogously, though the rate improvements over passive learning are more subtle.

6. Conclusions. Under Tsybakov’s noise conditions, active learning can offer improved asymptotic convergence rates compared to passive learning when the disagreement coefficient is finite. It is also possible to preserve these improved convergence rates when learning with a nested structure of hypothesis classes, using an algorithm that adapts to both the noise conditions and the complexity of the optimal classifier.

APPENDIX: DEFINITION OF $\hat{\mathcal{E}}$ AND RELATED QUANTITIES

We define the quantity $\hat{\mathcal{E}}_{\mathbb{C}}$ following Koltchinskii’s analysis of excess risk in terms of local Rademacher complexity [23]. The general idea is to construct a bound on the excess risk achieved by a given algorithm, such as empirical risk minimization, via an application of Talagrand’s inequality. Such a bound should be based on a measure of the expressiveness of the set of functions \mathbb{C} ; however, to bound the excess risk achieved by a particular algorithm given a number of data points, we need only measure the expressiveness of the set of functions the algorithm is likely to select from. For reasonable algorithms, such as empirical risk minimization, this means the set of functions with reasonably small excess risk. Thus, we can bound the excess risk of the algorithm in terms of a measure of expressiveness of the set of functions with relatively small risk, typically referred to as a *local* complexity measure. This reasoning is somewhat circular, in that first we must decide how small to expect the excess risk of the returned function to be before we can calculate the local complexity measure, which itself is used to calculate a bound on the risk of the returned function. Thus, we define the bound on the excess risk as a kind of fixed point. Furthermore, we can estimate these quantities using data-dependent confidence bounds, so that the excess risk bound can be calculated without direct access to the distribution. For the data-dependent measure of the expressiveness of the function class, we can use a Rademacher process. A detailed motivation and derivation can be found in [23].

For our purposes, we add an additional constraint, by requiring the functions we calculate the complexity of to agree with the labels of a labeled set \mathcal{L} . This is helpful for us, since given a set Q of labeled data with true labels, for any two functions h_1 and h_2 that agree on the labels of \mathcal{L} , it is always true that $er_{\mathcal{L} \cup Q}(h_1) - er_{\mathcal{L} \cup Q}(h_2)$ equals the difference of the true

empirical error rates. As we prove in the supplement, as long as the set \mathcal{L} is chosen carefully (i.e., as in Algorithm 2), the addition of this constraint is essentially inconsequential, so that $\hat{\mathcal{E}}_{\mathbb{C}}$ remains a valid excess risk bound. The detailed definitions are stated as follows.

For any function $f: \mathcal{X} \rightarrow \mathbb{R}$, and ξ_1, ξ_2, \dots a sequence of independent random variables with distribution uniform in $\{-1, +1\}$, define the *Rademacher process* for f under a finite set of (index, label) pairs $S \subset \mathbb{N} \times \mathcal{Y}$ as

$$R(f; S) = \frac{1}{|S|} \sum_{(i,y) \in S} \xi_i f(X_i).$$

The ξ_i should be thought of as internal variables in the learning algorithm, rather than being fundamental to the learning problem.

For any two finite sets $\mathcal{L} \subset \mathbb{N} \times \mathcal{Y}$ and $S \subset \mathbb{N} \times \mathcal{Y}$, define

$$\begin{aligned} \mathbb{C}[\mathcal{L}] &= \{h \in \mathbb{C} : er_{\mathcal{L}}(h) = 0\}, \\ \hat{\mathbb{C}}(\varepsilon; \mathcal{L}, S) &= \left\{h \in \mathbb{C}[\mathcal{L}] : er_S(h) - \min_{h' \in \mathbb{C}[\mathcal{L}]} er_S(h') \leq \varepsilon\right\}, \\ \hat{D}_{\mathbb{C}}(\varepsilon; \mathcal{L}, S) &= \sup_{h_1, h_2 \in \hat{\mathbb{C}}(\varepsilon; \mathcal{L}, S)} \frac{1}{|S|} \sum_{(i,y) \in S} \mathbb{1}[h_1(X_i) \neq h_2(X_i)] \end{aligned}$$

and

$$\hat{\phi}_{\mathbb{C}}(\varepsilon; \mathcal{L}, S) = \frac{1}{2} \sup_{h_1, h_2 \in \hat{\mathbb{C}}(\varepsilon; \mathcal{L}, S)} R(h_1 - h_2; S).$$

For $\delta, \varepsilon > 0$, $m \in \mathbb{N}$, define $s_m(\delta) = \ln \frac{20m^2 \log_2(3m)}{\delta}$ and $\mathbb{Z}_{\varepsilon} = \{j \in \mathbb{Z} : 2^j \geq \varepsilon\}$, and for any set $S \subset \mathbb{N} \times \mathcal{Y}$, define the set $S^{(m)} = \{(i, y) \in S : i \leq m\}$. We use the following definitions from Koltchinskii [23] with only minor modifications.

DEFINITION 3. For $\varepsilon \in [0, 1]$, and finite sets $S, \mathcal{L} \subset \mathbb{N} \times \mathcal{Y}$, define

$$\hat{U}_{\mathbb{C}}(\varepsilon, \delta; \mathcal{L}, S) = \hat{K} \left(\hat{\phi}_{\mathbb{C}}(\hat{c}\varepsilon; \mathcal{L}, S) + \sqrt{\frac{s_{|S|}(\delta) \hat{D}_{\mathbb{C}}(\hat{c}\varepsilon; \mathcal{L}, S)}{|S|}} + \frac{s_{|S|}(\delta)}{|S|} \right)$$

and

$$\hat{\mathcal{E}}_{\mathbb{C}}(S, \delta; \mathcal{L}) = \inf \left\{ \varepsilon > 0 : \forall j \in \mathbb{Z}_{\varepsilon}, \min_{m \in \mathbb{N}} \hat{U}_{\mathbb{C}}(2^j, \delta; \mathcal{L}^{(m)}, S^{(m)}) \leq 2^{j-4} \right\},$$

where, for our purposes, we can take $\hat{K} = 752$ and $\hat{c} = 3/2$, though there seems to be room for improvement in these constants. For completeness, we also define $\hat{\mathcal{E}}_{\mathbb{C}}(\emptyset, \delta; \mathbb{C}, \mathcal{L}) = \infty$ by convention.

We will also define a related quantity, representing a distribution-dependent version of $\hat{\mathcal{E}}$, also explored by Koltchinskii [23]. Specifically, for $\varepsilon > 0$, define

$$\mathbb{C}(\varepsilon) = \{h \in \mathbb{C} : er(h) - \nu \leq \varepsilon\}.$$

For $m \in \mathbb{N}$, let

$$\begin{aligned} \phi_{\mathbb{C}}(m, \varepsilon) &= \mathbb{E} \sup_{h_1, h_2 \in \mathbb{C}(\varepsilon)} |(er(h_1) - er_m(h_1)) - (er(h_2) - er_m(h_2))|, \\ \tilde{U}_{\mathbb{C}}(m, \varepsilon, \delta) &= \tilde{K} \left(\phi_{\mathbb{C}}(m, \tilde{c}\varepsilon) + \sqrt{\frac{s_m(\delta) \text{diam}(\tilde{c}\varepsilon; \mathbb{C})}{m}} + \frac{s_m(\delta)}{m} \right) \end{aligned}$$

and

$$\tilde{\mathcal{E}}_{\mathbb{C}}(m, \delta) = \inf\{\varepsilon > 0 : \forall j \in \mathbb{Z}_{\varepsilon}, \tilde{U}_{\mathbb{C}}(m, 2^j, \delta) \leq 2^{j-4}\},$$

where, for our purposes, we can take $\tilde{K} = 8272$ and $\tilde{c} = 3$. For completeness, we also define $\tilde{\mathcal{E}}_{\mathbb{C}}(0, \delta) = \infty$.

A.1. Definition of r_0 . In Definition 1, we took $r_0 = 0$. If $\theta < \infty$, then this choice is usually relatively harmless. However, in some cases, setting $r_0 = 0$ results in a suboptimal, or even infinite, value of θ , which is undesirable. In these cases, we would like to set r_0 as large as possible while maintaining the validity of the bounds. If we do this carefully enough, we should be able to establish bounds that, even in the worst case when $\theta = 1/r_0$, are never worse than the bounds for some analogous passive learning method; however, to do this requires r_0 to depend on the parameters of the learning problem: namely, n , δ , \mathbb{C} and \mathcal{D}_{XY} . The effect of a larger r_0 can sometimes be dramatic, as there are scenarios where $1 \ll \theta \ll 1/r_0$ [8]; we certainly wish to distinguish between such scenarios, and those where $\theta \propto 1/r_0$.

Generally, depending on the bound we wish to prove, different values of r_0 may be appropriate. For the tightest bound in terms of θ proven in the Appendices (namely, Lemma 7 of Appendix B in the supplementary material [20]), the definition of $r_0 = r_{\mathbb{C}}(n, \delta)$ in (13) below gives a good bound. For the looser bounds (namely, Theorems 5 and 6), a larger value of r_0 may provide better bounds; however, this same general technique can be employed to define a good value for r_0 in these looser bounds as well, simply using upper bounds on (13) analogous to how the theorems themselves are derived from Lemma 7 in Appendix B [20]. Likewise, one can state analogous refinements of r_0 for Theorems 1–4, though for brevity these are left for the reader's independent consideration.

DEFINITION 4. Define

$$(12) \quad \tilde{m}_{\mathbb{C}}(n, \delta) = \min \left\{ m \in \mathbb{N} : n \leq \log_2 \frac{4m^2}{\delta} + 2e \sum_{\ell=0}^{m-1} \mathbb{P}(\text{DIS}(\mathbb{C}(6\tilde{\mathcal{E}}_{\mathbb{C}}(\ell, \delta)))) \right\}$$

and

$$(13) \quad r_{\mathbb{C}}(n, \delta) = \max \left\{ \frac{1}{\tilde{m}_{\mathbb{C}}(n, \delta)} \sum_{\ell=0}^{\tilde{m}_{\mathbb{C}}(n, \delta)-1} \text{diam}(6\tilde{\mathcal{E}}_{\mathbb{C}}(\ell, \delta); \mathbb{C}), 2^{-n} \right\}.$$

We use this definition of $r_0 = r_{\mathbb{C}}(n, \delta)$ in all of the main proofs. In particular, with this definition, Lemma 7 of Appendix B [20] is never significantly worse than the analogous known result for passive learning (though it can be significantly better when $\theta \ll 1/r_0$).

Acknowledgments. I extend my sincere thanks to Larry Wasserman for numerous helpful discussions and also to John Langford for initially pointing out to me the possibility of A^2 adapting to Tsybakov’s noise conditions for threshold classifiers. I would also like to thank the anonymous referees for extremely helpful suggestions on earlier drafts.

SUPPLEMENTARY MATERIAL

Proofs and Supplements for “Rates of Convergence in Active Learning” (DOI: [10.1214/10-AOS843SUPP](https://doi.org/10.1214/10-AOS843SUPP); .pdf). The supplementary material contains three additional Appendices, namely, Appendices B, C and D. Specifically, Appendix B provides detailed proofs of Theorems 5–9, as well as several abstract lemmas from which these results are derived. Appendix C discusses the use of estimators in Algorithm 1. Finally, Appendix D includes a proof of a general minimax lower bound $\propto n^{-\kappa/(2\kappa-2)}$ for any nontrivial hypothesis class, generalizing a result of Castro and Nowak [12].

REFERENCES

- [1] ALEXANDER, K. S. (1984). Probability inequalities for empirical processes and a law of the iterated logarithm. *Ann. Probab.* **12** 1041–1067. [MR0757769](#)
- [2] ALEXANDER, K. S. (1985). Rates of growth for weighted empirical processes. In *Proceedings of the Berkeley Conference in Honor of Jerzy Neyman and Jack Kiefer II* 475–493. Wadsworth, Belmont, CA. [MR0822047](#)
- [3] ALEXANDER, K. S. (1986). Sample moduli for set-indexed gaussian processes. *Ann. Probab.* **14** 598–611. [MR0832026](#)
- [4] ALEXANDER, K. S. (1987). Rates of growth and sample moduli for weighted empirical processes indexed by sets. *Probab. Theory Related Fields* **75** 379–423. [MR0890285](#)
- [5] ANTHONY, M. and BARTLETT, P. L. (1999). *Neural Network Learning: Theoretical Foundations*. Cambridge Univ. Press, Cambridge. [MR1741038](#)
- [6] BALCAN, M.-F., BEYGEZIMER, A. and LANGFORD, J. (2006). Agnostic active learning. In *Proceedings of the 23rd International Conference on Machine Learning*. ACM, New York.
- [7] BALCAN, M.-F., BRODER, A. and ZHANG, T. (2007). Margin based active learning. In *Proceedings of the 20th Conference on Learning Theory. Lecture Notes in Computer Science* **4539** 35–50. Springer, Berlin. [MR2397577](#)

- [8] BALCAN, M.-F., HANNEKE, S. and WORTMAN, J. (2008). The true sample complexity of active learning. In *Proceedings of the 21st Conference on Learning Theory*. Omnipress, Madison, WI.
- [9] BALCAN, M.-F., BEYGEZIMIR, A. and LANGFORD, J. (2009). Agnostic active learning. *J. Comput. System Sci.* **75** 78–89. [MR2472318](#)
- [10] BEYGEZIMIR, A., DASGUPTA, S. and LANGFORD, J. (2009). Importance weighted active learning. In *International Conference on Machine Learning*. ACM, New York.
- [11] BLUMER, A., EHRENFUCHT, A., HAUSSLER, D. and WARMUTH, M. (1989). Learnability and the Vapnik-Chervonenkis dimension. *J. Assoc. Comput. Mach.* **36** 929–965. [MR1072253](#)
- [12] CASTRO, R. and NOWAK, R. (2008). Minimax bounds for active learning. *IEEE Trans. Inform. Theory* **54** 2339–2353. [MR2450865](#)
- [13] COHN, D., ATLAS, L. and LADNER, R. (1994). Improving generalization with active learning. *Machine Learning* **15** 201–221.
- [14] DASGUPTA, S., HSU, D. and MONTELEONI, C. (2007). A general agnostic active learning algorithm. In *Advances in Neural Information Processing Systems*. MIT Press, Cambridge, MA.
- [15] DEVROYE, L., GYÖRFI, L. and LUGOSI, G. (1996). *A Probabilistic Theory of Pattern Recognition*. Springer, New York. [MR1383093](#)
- [16] FRIEDMAN, E. (2009). Active learning for smooth problems. In *Proceedings of the 22nd Conference on Learning Theory*. Montreal, Quebec.
- [17] GINÉ, E. and KOLTCHINSKII, V. (2006). Concentration inequalities and asymptotic results for ratio type empirical processes. *Ann. Probab.* **34** 1143–1216. [MR2243881](#)
- [18] GINÉ, E., KOLTCHINSKII, V. and WELLNER, J. (2003). Ratio limit theorems for empirical processes. In *Stochastic Inequalities* (E. Giné, C. Houdré and D. Nualrt, eds.) 249–278. Birkhäuser, Basel. [MR2073436](#)
- [19] HANNEKE, S. (2007). A bound on the label complexity of agnostic active learning. In *Proceedings of the 24th International Conference on Machine Learning* (Z. Ghahramani, ed.) 353–360. ACM, New York.
- [20] HANNEKE, S. (2010). Proofs and supplements to “Rates of convergence in active learning.” DOI: [10.1214/10-AOS843SUPP](#).
- [21] KÄÄRIÄINEN, M. (2006). Active learning in the non-realizable case. In *Proceedings of the 17th International Conference on Algorithmic Learning Theory*. Springer, Berlin.
- [22] KEARNS, M. J., SCHAPIRE, R. E. and SELLIE, L. M. (1994). Toward efficient agnostic learning. *Mach. Learn.* **17** 115–141.
- [23] KOLTCHINSKII, V. (2006). Local rademacher complexities and oracle inequalities in risk minimization. *Ann. Statist.* **34** 2593–2656. [MR2329442](#)
- [24] KOLTCHINSKII, V. (2008). Oracle inequalities in empirical risk minimization and sparse recovery problems: Lecture notes. Technical report, Ecole d’été de Probabilités de Saint-Flour.
- [25] LI, Y. and LONG, P. M. (2007). Learnability and the doubling dimension. In *Advances in Neural Information Processing*. MIT Press, Cambridge, MA.
- [26] MAMMEN, E. and TSYBAKOV, A. (1999). Smooth discrimination analysis. *Ann. Statist.* **27** 1808–1829. [MR1765618](#)
- [27] MASSART, P. and NÉDÉLEC, E. (2006). Risk bounds for statistical learning. *Ann. Statist.* **34** 2326–2366. [MR2291502](#)
- [28] TSYBAKOV, A. B. (2004). Optimal aggregation of classifiers in statistical learning. *Ann. Statist.* **32** 135–166. [MR2051002](#)

- [29] VAN DER VAART, A. W. and WELLNER, J. A. (1996). *Weak Convergence and Empirical Processes*. Springer, New York. [MR1385671](#)
- [30] VAPNIK, V. (1982). *Estimation of Dependencies Based on Empirical Data*. Springer, New York. Translated from the Russian by Samuel Kotz. [MR0672244](#)
- [31] VAPNIK, V. (1998). *Statistical Learning Theory*. Wiley, New York. [MR1641250](#)
- [32] VAPNIK, V. and CHERVONENKIS, A. (1971). On the uniform convergence of relative frequencies of events to their probabilities. *Theory Probab. Appl.* **16** 264–280.
- [33] WANG, L. (2009). Sufficient conditions for agnostic active learnable. In *Advances in Neural Information Processing Systems 22*. MIT Press, Cambridge, MA.

DEPARTMENT OF STATISTICS
CARNEGIE MELLON UNIVERSITY
5000 FORBES AVENUE
PITTSBURGH, PENNSYLVANIA 15213
USA
E-MAIL: shanneke@stat.cmu.edu